



## Query-Oriented Summarization of RDF Graphs

Šejla Čebirić, François Goasdoué, Ioana Manolescu

### ► To cite this version:

Šejla Čebirić, François Goasdoué, Ioana Manolescu. Query-Oriented Summarization of RDF Graphs. Proceedings of the VLDB Endowment, Aug 2015, Kohala Coast, Hawaii, United States. hal-01178140

**HAL Id: hal-01178140**

**<https://hal.inria.fr/hal-01178140>**

Submitted on 17 Jul 2015

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Query-Oriented Summarization of RDF Graphs

Šejla Čebirić  
INRIA & U. Paris-Sud, France  
sejla.cebirc@inria.fr

François Goasdoué  
U. Rennes 1 & INRIA, France  
fg@irisa.fr

Ioana Manolescu  
INRIA & U. Paris-Sud, France  
ioana.manolescu@inria.fr

## 1. INTRODUCTION

The Resource Description Framework (RDF) is a graph-based data model promoted by the W3C as the standard for Semantic Web applications. Its associated query language is SPARQL. RDF graphs are often *large* and *varied*, produced in a variety of contexts, e.g., scientific applications, social or online media, government data etc. They are *heterogeneous*, i.e., resources described in an RDF graph may have very different sets of properties. An RDF resource may have: no types, one or several types (which may or may not be related to each other). *RDF Schema* (RDFS) information may optionally be attached to an RDF graph, to enhance the description of its resources. Such statements also entail that in an RDF graph, some data is **implicit**. According to the W3C RDF and SPARQL specification, **the semantics of an RDF graph comprises both its explicit and implicit data**; in particular, SPARQL query answers must be computed *reflecting both the explicit and implicit data*. These features make RDF graphs complex, both structurally and conceptually. It is intrinsically hard to get familiar with a new RDF dataset, especially if an RDF schema is sparse or not available at all.

In this work, we study the problem of *RDF summarization*, that is: given an input RDF graph  $G$ , find an RDF graph  $S_G$  which *summarizes  $G$  as accurately as possible, while being possibly orders of magnitude smaller* than the original graph. Such a summary can be used in a variety of contexts: to help an RDF application designer get acquainted with a new dataset, as a first-level user interface, or as a support for query optimization as typically used in semi-structured graph data management [4] etc. Our approach is *query-oriented*, i.e., a summary should enable static analysis and help formulating and optimizing queries; for instance, querying a summary of a graph should reflect whether the query has some answers against this graph, or finding a simpler way to formulate the query etc. Ours is the first semi-structured data summarization approach focused on *partially explicit, partially implicit* RDF graphs.

In the sequel, Section 2 recalls RDF basics, and sets the requirements for our query-oriented RDF summaries. Section 3 describes two flavors of summaries: a *baseline* which is compact, simple and meets our requirements, at the expense of a strong simplification of the graph, and a *refined* one which trades some of our requirements for more accuracy in representing the structure. Section 4 presents our scenario, we then discuss related work and conclude.

Assertion	Triple	Relational notation
Class	$s \text{ rdf:type } o$	$o(s)$
Property	$s \text{ p } o$	$p(s, o)$

Constraint	Triple	OWA interpretation
Subclass	$s \text{ rdfs:subClassOf } o$	$s \subseteq o$
Subproperty	$s \text{ rdfs:subPropertyOf } o$	$s \subseteq o$
Domain typing	$s \text{ rdfs:domain } o$	$\Pi_{\text{domain}}(s) \subseteq o$
Range typing	$s \text{ rdfs:range } o$	$\Pi_{\text{range}}(s) \subseteq o$

Figure 1: RDF (top) & RDFS (bottom) statements.

## 2. PRELIMINARIES

We introduce RDF graph and queries in Section 2.1, and requirements for our RDF summaries in Section 2.2.

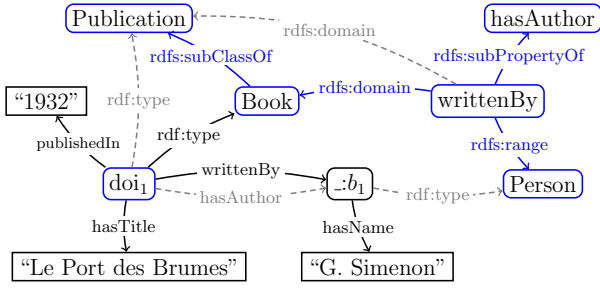
### 2.1 RDF Graphs and Queries

An *RDF graph* (or *graph*, in short) is a set of *triples* of the form  $s \text{ p } o$ , stating that the *subject*  $s$  has the *property*  $p$ , and the value of that property is the *object*  $o$ .

We consider only well-formed triples, as per the W3C's RDF specification, using uniform resource identifiers (URIs), typed or untyped literals (constants), and *blank nodes* (unknown URIs or literals) corresponding to a form of incomplete information, similar to *unknown URI* or *literal tokens*. **Notations** We use  $s$ ,  $p$ , and  $o$  in triples as placeholders. Literals are shown as strings between quotes, e.g., “*string*”.

Figure 1 (top) shows how to use triples to describe resources: resource types are described through unary relations, properties using binary relations. The RDF standard provides a set of built-in classes and properties in the `rdf:` and `rdfs:` pre-defined namespaces, e.g., triples of the form  $s \text{ rdf:type } o$  specify the class(es) to which a resource belongs. *For brevity, we will sometimes use  $\tau$  to denote `rdf:type`*. For example, the RDF graph  $G$  shown in Figure 2 describes a book, identified by `doi1`: its author (a blank node `:b1` related to the author name), title and publication date.

**RDF Schema** allows enhancing the descriptions in RDF graphs by means of *RDFS triples*, declaring *semantic constraints* between the graph classes and properties. The RDFS constraints (Figure 1) lead to **implicit triples** which may be part of an RDF graph even though they are not physically present in it. An implicit triple can be obtained by an *immediate entailment step* based on (i) an RDFS constraint, and (ii) either a second constraint (also called *schema triple*) or an RDF triple that is not a constraint (also termed *data triple*). A triple is *entailed by a graph  $G$* , if and only if there is a sequence of applications of entailment rules that leads from the graph to the triple (where at each step of the entail-



$G = \{ \text{doi}_1 \text{ rdfs:type Book, doi}_1 \text{ writtenBy } \text{:b}_1, \text{doi}_1 \text{ hasTitle "Port des Brumes", } \text{:b}_1 \text{ hasName "G. Simenon", doi}_1 \text{ publishedIn "1932"} \}$

Figure 2: Sample RDF graph.

ment sequence, the triples previously entailed are also taken into account). For instance, assume that the RDF graph  $G$  above is extended with the following constraints:

- **Book** `rdfs:subClassOf` **Publication**
- **writtenBy** `rdfs:subPropertyOf` **hasAuthor**
- **writtenBy** `rdfs:domain` **Book** and **writtenBy** `rdfs:range` **Person**

The resulting graph is depicted in Figure 2. Its implicit triples are those represented by dashed-line edges.

**Saturation** The immediate entailment rules allow defining the finite *saturation* (a.k.a. closure) of an RDF graph  $G$ , which is the RDF graph  $G^\infty$  defined as the fixed-point obtained by repeatedly applying entailment rules on  $G$ .

The saturation of an RDF graph is unique (up to blank node renaming), and does not contain implicit triples (they have all been made explicit by saturation). Clearly, a graph  $G$  entails a triple if and only if the triple belongs to  $G^\infty$ . RDF entailment is part of the RDF standard; *the answers to a query posed on  $G$  must take into account all triples in  $G^\infty$ , since the semantics of an RDF graph is its saturation.*

**Queries** We consider the SPARQL dialect consisting of *basic graph pattern* (BGP) queries, a.k.a. conjunctive queries, widely considered in research but also in real-world applications [8]. A BGP is a set of *query triple patterns*, or query triples in short; each triple has a subject, property and object, some of which can be variables.

**Query answering** The evaluation of a query  $q$  against  $G$  has access only to  $G$ 's explicit triples, thus may lead to an incomplete answer; the complete answer is obtained by evaluating  $q$  against  $G^\infty$ . For instance, the query below asks for name of the author of "Le Pont des Brumes":

$q(x_3) :- x_1 \text{ hasAuthor } x_2, x_2 \text{ hasName } x_3$   
 $x_1 \text{ hasTitle "Le Port des Brumes"}$

Its answer against the graph in Figure 2 is  $q(G^\infty) = \{ \text{"G. Simenon"} \}$ . Note that evaluating  $q$  only against  $G$  leads to the empty answer, which is obviously incomplete.

## 2.2 RDF Summary Requirements

We assume that *the summary  $S_G$  of an RDF graph  $G$  is an RDF graph itself*. Further, we require the following:

**Completeness** The saturation of the summary of  $G$  must be the same as the summary of its saturation  $G^\infty$ , due to the semantics of an RDF graph being its saturation.

**Schema independence** It must be possible to summarize  $G$  whether or not it has associated RDFS triples.

The following properties are of a more quantitative nature:

**Compactness** The summary should be typically smaller than the RDF graph, ideally by orders of magnitude.

**Representativeness** The summary should not lose too much information from  $G$ .

**Accuracy** The summary should avoid, to the extent possible, reflecting data that does not exist in  $G$ .

A trade-off exists between compactness and representativeness, as the latter tends to require more information.

**Criteria for representativeness and accuracy** Our query-oriented RDF graph summarization leads us to the following criteria. For *representativeness*, queries with results on  $G$  should also have results on the summary. Symmetrically, for *accuracy*, a query that can be matched on the summary, should also be matched on the RDF graph itself. To formalize these, let  $\mathcal{Q}$  be a SPARQL dialect.

**DEFINITION 1. (QUERY-BASED REPRESENTATIVENESS)**  $S_G$  is  $\mathcal{Q}$ -representative of  $G$  if and only if for any query  $q \in \mathcal{Q}$  such that  $q(G^\infty) \neq \emptyset$ , we have  $q(S_G^\infty) \neq \emptyset$ .

Note that *several graphs may have the same summary*, since a summary loses *some* of the information from the original graph. If two RDF graphs differ only with respect to such information, they have the same summary. We term *inverse set* of a summary  $S_G$ , the set of all RDF graphs whose summary is  $S_G$ . This leads to the accuracy criterion, with respect to *any graph a summary may correspond to*:

**DEFINITION 2. (QUERY-BASED ACCURACY)** Let  $S_G$  be a summary, and  $\mathcal{G}$  the inverse set of  $S_G$ . The summary  $S_G$  is  $\mathcal{Q}$ -accurate if for any query  $q \in \mathcal{Q}$  such that  $q(S_G^\infty) \neq \emptyset$ , there exists  $G \in \mathcal{G}$  such that  $q(G^\infty) \neq \emptyset$ .

For compactness, the (voluminous) set of literals, along with subject and object URIs for non- $\tau$  triples from  $G$  should not appear in  $S_G$ . However, given that property URIs are often specified in SPARQL queries, and that typically there are far less distinct property URIs than the subject or object URIs [11], property URIs should be preserved by the summary. This leads us to the following SPARQL dialect:

**DEFINITION 3. (RELATIONAL BGP)** A relational BGP (RBGP, in short) is a BGP query whose body has: (i) URIs in all the property positions, (ii) a URI in the object position of every  $\tau$  triple, and (iii) variables in any other positions.

We define *RBGP representativeness* and *RBGP accuracy* by instantiating  $\mathcal{Q}$  in Definition 1 and Definition 2, respectively, to RBGP queries (Definition 3).

## 3. RDF SUMMARIES

This section describes our RDF summary proposals.

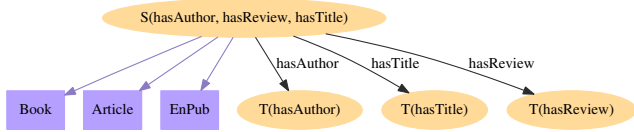


Figure 3: Baseline summary  $B_{G_0}$  for the sample graph  $G_0$ .

### 3.1 Baseline Summary

We assume a function `newURI()` returning a fresh URI on each call. We call *data property* any property  $p$  occurring in the data component of  $G$ , and different from  $\tau$ . Further, for any data property  $p$ , the *property source* of  $p$ , denoted  $S(p)$ , is a URI set by `newURI()`, and similarly, the *property target* of  $p$ , denoted  $T(p)$ , is a URI set by `newURI()`.

**DEFINITION 4. (BASELINE SUMMARY)** *Given an RDF graph  $G$ , the baseline summary of  $G$  is an RDF graph  $B_G$  such that:*

**Schema**  $B_G$  has the same schema triples as  $G$ .

**DNT** (Data triples of  $B_G$  whose property is not  $\tau$ ) Let  $p, p_1, p_2$  be some data properties from  $G$ .

**DNT1** The triple  $S(p) \ p \ T(p)$  belongs to  $B_G$ ;

**DNT2** if  $s \ p_1 \ o_1, s \ p_2 \ o_2 \in G$ , then  $S(p_1) = S(p_2)$ ;

**DNT3** if  $s_1 \ p_1 \ o, s_2 \ p_2 \ o \in G$ , then  $T(p_1) = T(p_2)$ ;

**DNT4** if  $s \ p_1 \ o_1, o_1 \ p_2 \ o_2 \in G$ , then  $T(p_1) = S(p_2)$ ;

**DT** (Data triples of  $B_G$  whose property is  $\tau$ )

**DT1** If  $s \ p \ o, s \ \tau \ c$  are in  $G$ , then  $S(p) \ \tau \ c$  is in  $B_G$ ;

**DT2** if  $s \ p \ o, o \ \tau \ c$  are in  $G$ , then  $T(p) \ \tau \ c$  is in  $B_G$ ;

**DT3** Let  $n_\tau$  be set to `newURI()`. If  $s \ \tau \ c \in G$ , and  $\nexists s \ p \ o \in G$  and  $\nexists s' \ p \ s \in G$ , then  $n_\tau \ \tau \ c \in B_G$ .

The baseline summary has the same schema as  $G$  (**Schema**), as well as a source and a target URI for each data property (**DNT1**). As soon as two data properties have the same subject and/or object, their corresponding source and target URIs are the same in  $B_G$  accordingly (**DNT2-DNT4**). If the subject (or the object) of a data property  $p$  is of type  $c$ , then the  $p$  source (or target) URI is declared to be of type  $c$  in  $B_G$  (**DT1-DT2**). Finally, a single URI in  $B_G$  reflects all the subjects of  $\tau$  triples in  $G$  that do not appear as subject or object of a data property in  $G$  (**DT3**).

For instance, consider the graph  $G_0$  describing the resources *rdoi1* to *rdoi3*:

<i>rdoi1</i> $\tau$ <i>rBook</i>	<i>rdoi2</i> $\tau$ <i>rEnPub</i>
<i>rdoi1</i> <i>rhasTitle</i> <i>r</i> "T1"	<i>rdoi2</i> $\tau$ <i>rArticle</i>
<i>rdoi1</i> <i>rhasAuthor</i> <i>r</i> "A1"	<i>rdoi2</i> <i>rhasTitle</i> <i>r</i> "T2"
<i>rdoi1</i> <i>rhasAuthor</i> <i>r</i> "A2"	<i>rdoi2</i> <i>rhasAuthor</i> <i>r</i> "A3"
<i>rdoi3</i> <i>rhasTitle</i> <i>r</i> "T3"	<i>rdoi2</i> <i>rhasReview</i> <i>r</i> "R1"
<i>rdoi3</i> <i>rhasAuthor</i> <i>r</i> "A4"	

where *rEnPub* is the class of publications in English; we omitted a schema, to focus on the treatment of the data triples. Figure 3 depicts the baseline summary of  $G_0$ ; the rectangular nodes correspond to class URIs copied from  $G_0$ , whereas the oval nodes are URIs created by `newURI()`.

Importantly, the baseline summary meets our requirements, as follows (the proofs can be found in [12]). We say two summary graphs are *equivalent*, denoted  $\equiv$ , iff they are

identical up to a bijection between their sets of URIs. The completeness requirement is met with Proposition 1, stating the commutativity of saturation and summarization:

**PROPOSITION 1.** *Let  $B_G$  be the baseline summary of  $G$ , and  $B_{G^\infty}$  the baseline summary of  $G^\infty$ . Then:  $(B_G)^\infty \equiv B_{G^\infty}$ .*

Regarding representativeness and accuracy, we show [12]:

**PROPOSITION 2.** *The baseline summary is (i) RBGP representative and (ii) RBGP accurate.*

It is easy to see that the baseline summary can be built in  $O(|G|^2)$  time. Its size is bounded by the size of  $G$ 's schema to which we add (i) the number of data properties from  $G$  and (ii) the number of class assertions from  $G$ .

### 3.2 Refined Summary

The baseline summary may unify property source and target URIs quite aggressively. For instance, if a *store* and a *person* both have a *zipcode*, they will lead to the same baseline URI (through rule **DNT2**), even though they are very different things.

To mitigate this issue, we designed a second flavor of summary of an RDF graph  $G$ , termed *refined* and denoted  $R_G$ . For space reasons, the definition is delegated to [12]. Intuitively, the difference between the baseline and the refined summary is that the latter fuses data property source and/or target URIs *only if one resource in  $G$  that leads to their unification has no type at all*. For illustration, the refined summary  $R_{G_0}$  of the same sample graph appears in Figure 4. The target URIs  $T(rhasTitle)$ ,  $T(rhasAuthor)$  and  $T(rhasReview)$  are the same as in  $B_{G_0}$ ; however, in  $R_{G_0}$ , four URIs have been created in the upper row to represent resources having both a title and an author, respectively (from left to right): resources of type *rBook*, those having no type in  $G_0$ , those of type *rArticle*, and those of type *rEnPub*.

The refined summary commutes with saturation (given a graph  $G$ ,  $(R_G)^\infty = R_{G^\infty}$ ); it is also RBGP accurate [12]. It is *more accurate than the baseline*, as illustrated in Figure 4: the rightmost URI in  $R_{G_0}$  shows that only resources of type *rArticle* or *rEnPub* have reviews, whereas the baseline  $B_{G_0}$  may lead one to believe that a resource of type *rBook* also has a review. (Recall that this may happen in *some* graph  $G_1$  whose summary is  $B_{G_0}$ ; it just does not happen in  $G_0$ ).

This extra accuracy comes at a cost. Computing the refined summary has  $O(|G|^5)$  complexity, which requires an efficient underlying system e.g., based on triple partitioning and indexing etc. The refined summary is representative for *all RBGPs which do not have more than one  $\tau$  triple with the same subject*. This follows from a *graph homomorphism* from  $G^\infty$  to  $(R_G)^\infty$  [12].

An upper bound for its size is the number of classes in  $G \times$  the number of distinct data properties. This is significantly larger than for  $B_G$ , but in practice (i) the bound is seldom reached (ii) more accurate summaries are better appreciated by users getting acquainted with RDF graphs.

## 4. DEMONSTRATION SCENARIO

We demonstrate our Java tool (7.700 lines approx.) for computing baseline and refined RDF graph summaries. The tool issues queries that are executed by the underlying RDF store, in particular OpenLink Virtuoso Server (7.1), a PostgreSQL based RDF store complete with indexes etc., and a Hadoop-based RDF query processing platform [2].

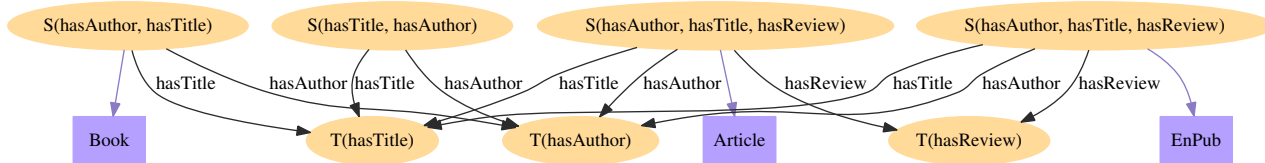


Figure 4: Refined summary  $R_{G_0}$  for the sample graph  $G_0$ .

Demo attendees will be able to: (i) pick an RDF graph  $G$  from a set including LUBM data, an RDF dump of DBLP, open data sets from the French INSEE (statistics) and IGN (geographic) institutes, as well as small hand-crafted examples chosen for their interest in illustrating summary features such as representativity, accuracy, dependence on the degree of saturation (recall that the saturated summary does not depend on whether  $G$  is saturated, but the non-saturated  $B_G$  and  $R_G$  do!) (ii) compute  $B_G$  and  $R_G$  using one of the systems; (iii) inspect the summary with the help of ATT’s GraphViz/DOT-based GUI; (iv) trigger the saturation of the summary; (v) modify the graphs in the store and see the impact on the summaries; (vi) choose or write custom RBGP queries, comprising *property paths*, specified by regular expressions of SPARQL v1.1., and see how they *unfold* into unions of queries, or how the summary allows deciding that they have empty results. Step (vi) adapts Dataguide techniques [4] to our RDF-specific summaries.

## 5. RELATED WORK

**OEM and XML summaries** Dataguides [4] were introduced to summarize semistructured OEM graphs, similar to RDF, but assumed to have a “root” node, from which all others are accessible; this may not hold for RDF. Dataguides construction has worst-case exponential time complexity, thus is not in general feasible. Many works considered indexes for supporting XML path queries; these works differ from ours, because the input is a tree or DAG and/or because XML lacks types and implicit information.

**Graph summarization** Graph summarization has been very intensively studied, in particular through mining or clustering; large-scale graph processing is also a hot topic. The notions of *summaries* and *structural indexes* bear similarities, both being a reduced version of the input graph due to collapsing nodes based on some common attributes. Our focus is on *RDF graphs with implicit data*, for which we devised *query-oriented summaries*, which are RDF graphs themselves and may be computed on a variety of platforms.

A graph core  $C$  for a given graph  $G$  is a graph such that an isomorphism exists between  $G$  and  $C$ , and  $C$  is the smallest graph with this property [3]. Neither of our summaries are cores of  $G$ , since a homomorphism is not guaranteed to exist from either summary to  $G$ . In exchange, both our summaries can be built in polynomial time in the size of  $G$ , while computing the core is much harder.

In [7], graphs from the LOD cloud are summarized, focusing on the distribution of classes and properties across LOD sources. The size of a bisimulation may explode exponentially w.r.t. that of  $G$ . As we aim for both complete and compact summaries bisimulation is not a good fit.

To overcome bisimulation issues, [5] suggests locality-based summaries, whose generation requires removing the (many) triples whose objects are literals from the input graph. Our summaries represent the queries over these triples as well.

A *triple-oriented* structural index for RDF data is built in [9] as a non-RDF graph. Each node represents a set of triples, while edges describe how the triples from adjacent nodes join. [10] proposes a tree RDF index, storing regions defined by center vertices, limited to property paths and built assuming that the input RDF graph is saturated.

In [6], RDF classes are inferred based on the common properties of resources. Thus, only common source patterns are analyzed, while the common targets and property paths are not considered; `rdf:type` triples are also ignored. [1] explores alternative RDF summaries w.r.t. graph homomorphism and trades precision for computing efficiency.

Finally, summarizing implicit data is not considered in any of these works.

**Acknowledgments** This work has been partially funded by the projects PIA Datalyse and DGA RAPID ODIN.

## 6. REFERENCES

- [1] Stéphane Campinas, Renaud Delbru, and Giovanni Tummarello. Efficiency and precision trade-offs in graph summary algorithms. In *IDEAS*, 2013.
- [2] François Goasdoué, Zoi Kaoudi, Ioana Manolescu, Jorge-Arnulfo Quiané-Ruiz, and Stamatis Zampetakis. CliqueSquare: Flat Plans for Massively Parallel RDF Queries. In *ICDE*, 2015.
- [3] Chris Godsil and Gordon Royle. *Algebraic graph theory*. Springer-Verlag, 2001.
- [4] Roy Goldman and Jennifer Widom. Dataguides: Enabling query formulation and optimization in semistructured databases. In *VLDB*, 1997.
- [5] Sairam Gurajada, Stephan Seufert, Iris Miliaraki, and Martin Theobald. Using graph summarization for join-ahead pruning in a distributed RDF engine. In *SWIM workshop*, 2014.
- [6] Kenza Kellou-Menouer and Zoubida Kedad. A clustering based approach for type discovery in RDF data sources. In *EGC*, 2015.
- [7] Shahan Khatchadourian and Mariano P. Consens. ExpLOD: Summary-based exploration of interlinking and RDF usage in the linked open data cloud. In *ESWC*, 2010.
- [8] Davide Lanti, Martin Rezk, Guohui Xiao, and Diego Calvanese. The NPD benchmark: Reality check for OBDA systems. In *EDBT*, 2015.
- [9] François Picalausa, Yongming Luo, George H. L. Fletcher, Jan Hidders, and Stijn Vansummeren. A structural approach to indexing triples. In *ESWC*, 2012.
- [10] Octavian Udrea, Andrea Pugliese, and V. S. Subrahmanian. GRIN: A graph based RDF index. In *AAAI*. AAAI Press, 2007.
- [11] Statistics on The Billion Triple Challenge Dataset. [gromgull.net/blog/2010/09/btc2010-basic-stats](http://gromgull.net/blog/2010/09/btc2010-basic-stats), 2010.
- [12] Extended version of this work (technical report), 2015.